

Efficient Estimation of Multidimensional Regression Model with Multilayer Perceptron

Joseph Rynkiewicz¹

Université Paris I - SAMOS/MATISSE
72 rue Regnault, Paris - France

Abstract. This work concerns estimation of multidimensional nonlinear regression models using multilayer perceptron (MLP). The main problem with such model is that we have to know the covariance matrix of the noise to get optimal estimator. However we show that, if we choose as cost function the logarithm of the determinant of the empirical error covariance matrix, we get an asymptotically optimal estimator.

1 Introduction

Let us consider a sequence $(Y_t, Z_t)_{t \in \mathbb{N}}$ of i.i.d.¹ random vectors (i.e. identically distributed and independent). So, each couple (Y_t, Z_t) has the same law that a generic variable (Y, Z) . Moreover, we assume that the model can be written

$$Y_t = F_{W^0}(Z_t) + \varepsilon_t$$

where

- F_{W^0} is a function represented by a MLP with parameters or weights W^0 .
- (ε_t) is an i.i.d. centered noise with unknown invertible covariance matrix Γ_0 .

Our goal is to estimate the true parameter by minimizing an appropriate cost function. This model is called a regression model and a popular choice for the associated cost function is the mean square error :

$$\frac{1}{n} \sum_{t=1}^n \|Y_t - F_W(Z_t)\|^2$$

where $\|\cdot\|$ denotes the Euclidean norm on \mathbb{R}^d . Although this function is widely used, it is easy to show that we get then a suboptimal estimator. An other solution is to use an approximation of the covariance error matrix to compute generalized least square estimator :

$$\frac{1}{n} \sum_{t=1}^n (Y_t - F_W(Z_t))^T \Gamma^{-1} (Y_t - F_W(Z_t)),$$

¹It is not hard to extend all what we show in this paper for stationary mixing variables and so for time series

where T denotes the transposition of the matrix. Here we assume that Γ is a good approximation of the true covariance matrix of the noise Γ_0 . However it takes time to compute a good approximation of matrix Γ_0 and it leads asymptotically to the cost function proposed in this article (see for example Rynkiewicz [4]) :

$$U_n(W) := \log \det \left(\frac{1}{n} \sum_{t=1}^n (Y_t - F_W(Z_t))(Y_t - F_W(Z_t))^T \right) \quad (1)$$

This paper is devoted to the theoretical study of $U_n(W)$. We assume that the true architecture of the MLP is known so that the Hessian matrix computed in the sequel verifies the assumption to be definite positive (see Fukumizu [1]).

In this framework, we study the asymptotic behavior $\hat{W}_n := \arg \min U_n(W)$, the weights minimizing the cost function $U_n(W)$. We show that under simple assumptions this estimator is asymptotically optimal in the sense that it has the same asymptotic behavior than the generalized least square estimator using the true covariance matrix of the noise.

Numerical procedures to compute this estimator and examples of its behavior can be found in Rynkiewicz [4].

2 The first and second derivatives of $W \mapsto U_n(W)$

First, we introduce a notation : if $F_W(X)$ is a d -dimensional parametric function depending of a parameter W , let us write $\frac{\partial F_W(X)}{\partial W_k}$ (resp. $\frac{\partial^2 F_W(X)}{\partial W_k \partial W_l}$) for the d -dimensional vector of partial derivative (resp. second order partial derivatives) of each component of $F_W(X)$.

2.1 First derivatives

Now, if $\Gamma_n(W)$ is a matrix depending of the parameter vector W , we get From Magnus and Neudecker [3]

$$\frac{\partial}{\partial W_k} \ln \det (\Gamma_n(W)) = \text{tr} \left(\Gamma_n^{-1}(W) \frac{\partial}{\partial W_k} \Gamma_n(W) \right)$$

here

$$\Gamma_n(W) = \frac{1}{n} \sum_{t=1}^n (y_t - F_W(z_t))(y_t - F_W(z_t))^T$$

note that these matrix $\Gamma_n(W)$ and its inverse are symmetric. Now, if we note

$$A_n(W_k) = \frac{1}{n} \sum_{t=1}^n \left(-\frac{\partial F_W(z_t)}{\partial W_k} (y_t - F_W(z_t))^T \right)$$

using the fact

$$\text{tr} (\Gamma_n^{-1}(W) A_n(W_k)) = \text{tr} (A_n^T(W_k) \Gamma_n^{-1}(W)) = \text{tr} (\Gamma_n^{-1}(W) A_n^T(W_k))$$

we get

$$\frac{\partial}{\partial W_k} \ln \det (\Gamma_n(W)) = 2 \text{tr} (\Gamma_n^{-1}(W) A_n(W_k)) \quad (2)$$

2.2 Second derivatives

We write now

$$B_n(W_k, W_l) := \frac{1}{n} \sum_{t=1}^n \left(\frac{\partial F_W(z_t)}{\partial W_k} \frac{\partial F_W(z_t)}{\partial W_l}^T \right)$$

and

$$C_n(W_k, W_l) := \frac{1}{n} \sum_{t=1}^n \left(-(y_t - F_W(z_t)) \frac{\partial^2 F_W(z_t)}{\partial W_k \partial W_l}^T \right)$$

We get

$$\begin{aligned} \frac{\partial^2 U_n(W)}{\partial W_k \partial W_l} &= \frac{\partial}{\partial W_l} 2tr \left(\Gamma_n^{-1}(W) A_n(W_k) \right) = \\ &2tr \left(\frac{\partial \Gamma_n^{-1}(W)}{\partial W_l} A_n(W_k) \right) + 2tr \left(\Gamma_n^{-1}(W) B_n(W_k, W_l) \right) + 2tr \left(\Gamma_n(W)^{-1} C_n(W_k, W_l) \right) \end{aligned}$$

Now, Magnus and Neudecker [3] give an analytic form of the derivative of an inverse matrix, so we get

$$\begin{aligned} \frac{\partial^2 U_n(W)}{\partial W_k \partial W_l} &= 2tr \left(\Gamma_n^{-1}(W) \left(A_n(W_k) + A_n^T(W_k) \right) \Gamma_n^{-1}(W) A_n(W_k) \right) + \\ &2tr \left(\Gamma_n^{-1}(W) B_n(W_k, W_l) \right) + 2tr \left(\Gamma_n^{-1}(W) C_n(W_k, W_l) \right) \end{aligned}$$

so

$$\begin{aligned} \frac{\partial^2 U_n(W)}{\partial W_k \partial W_l} &= 4tr \left(\Gamma_n^{-1}(W) A_n(W_k) \Gamma_n^{-1}(W) A_n(W_k) \right) \\ &+ 2tr \left(\Gamma_n^{-1}(W) B_n(W_k, W_l) \right) + 2tr \left(\Gamma_n^{-1}(W) C_n(W_k, W_l) \right) \end{aligned} \quad (3)$$

3 Asymptotic properties of \hat{W}_n

First, following the same lines that Yao [5], it is easy to show that, if the noise of the model has a moment of order at least 2, the estimator is strongly consistent (i.e. $\hat{W}_n \xrightarrow{a.s.} W^0$).

Moreover, for a MLP function, there exists a constant C such that we have the following inequalities :

$$\begin{aligned} \left\| \frac{\partial F_W(Z)}{\partial W_k} \right\| &\leq C(1 + \|Z\|) \\ \left\| \frac{\partial^2 F_W(Z)}{\partial W_k \partial W_l} \right\| &\leq C(1 + \|Z\|^2) \\ \left\| \frac{\partial^2 F_W(Z)}{\partial W_k \partial W_l} - \frac{\partial^2 F_W^0(Z)}{\partial W_k \partial W_l} \right\| &\leq C\|W - W^0\|(1 + \|Z\|^3) \end{aligned}$$

So, if Z has a moment of order at least 3 (see the justification in Yao [5]), we get the following lemma :

Lemma 1 *Let $\Delta U_n(W^0)$ be the gradient vector of $U_n(W)$ at W^0 , $\Delta U(W^0)$ be the gradient vector of $U(W) := \log \det(Y - F_W(Z))$ at W^0 and $HU_n(W^0)$ be the Hessian matrix of $U_n(W)$ at W^0 .*

We define finally

$$B(W_k, W_l) := \frac{\partial F_W(Z)}{\partial W_k} \frac{\partial F_W(Z)}{\partial W_l}^T$$

and

$$A(W_k) = \left(-\frac{\partial F_W(Z)}{\partial W_k} (Y - F_W(Z))^T \right)$$

We get then

1. $HU_n(W^0) \xrightarrow{a.s.} 2I_0$
2. $\sqrt{n}\Delta U_n(W^0) \xrightarrow{Law} \mathcal{N}(0, 4I_0)$

where, the component (k, l) of the matrix I_0 is :

$$tr \left(\Gamma_0^{-1} E \left(B(W_k^0, W_l^0) \right) \right)$$

proof To prove the lemma, we remark first that the component (k, l) of the matrix $4I_0$ is :

$$E \left(\frac{\partial U(W^0)}{\partial W_k} \frac{\partial U(W^0)}{\partial W_l} \right) = E \left(2tr \left(\Gamma_0^{-1} A^T(W_k^0) \right) \times 2tr \left(\Gamma_0^{-1} A(W_l^0) \right) \right)$$

and, since the trace of the product is invariant by circular permutation,

$$\begin{aligned} E \left(\frac{\partial U(W^0)}{\partial W_k} \frac{\partial U(W^0)}{\partial W_l} \right) &= \\ 4E \left(-\frac{\partial F_{W^0}(Z)}{\partial W_k} \Gamma_0^{-1} (Y - F_{W^0}(Z)) (Y - F_{W^0}(Z))^T \Gamma_0^{-1} \left(-\frac{\partial F_{W^0}(Z)}{\partial W_l} \right) \right) &= \\ 4E \left(\frac{\partial F_{W^0}(Z)}{\partial W_k} \Gamma_0^{-1} \frac{\partial F_{W^0}(Z)}{\partial W_l} \right) &= \\ 4tr \left(\Gamma_0^{-1} E \left(\frac{\partial F_{W^0}(Z)}{\partial W_k} \frac{\partial F_{W^0}(Z)}{\partial W_l} \right) \right) &= \\ 4tr \left(\Gamma_0^{-1} E \left(B(W_k^0, W_l^0) \right) \right) \end{aligned}$$

Now, for the component (k, l) of the expectation of the Hessian matrix, we remark that

$$\lim_{n \rightarrow \infty} tr \left(\Gamma_n^{-1}(W^0) A_n(W_k^0) \Gamma_n^{-1}(W^0) A_n(W_k^0) \right) = 0$$

and

$$\lim_{n \rightarrow \infty} tr \Gamma_n^{-1} C_n(W_k^0, W_l^0) = 0$$

so

$$\begin{aligned} \lim_{n \rightarrow \infty} H_n(W^0) &= \lim_{n \rightarrow \infty} 4tr \left(\Gamma_n^{-1}(W^0) A_n(W_k^0) \Gamma_n^{-1}(W^0) A_n(W_k^0) \right) + \\ 2tr \Gamma_n^{-1}(W^0) B_n(W_k^0, W_l^0) &+ 2tr \Gamma_n^{-1} C_n(W_k^0, W_l^0) = \\ = 2tr \left(\Gamma_0^{-1} E \left(B(W_k^0, W_l^0) \right) \right) \end{aligned}$$

■

From a classical argument of local asymptotic normality (see for example Yao [5]), we deduce then the following property for the estimator \hat{W}_n :

Proposition 1 *Let W_n^* the estimator of the generalized least square :*

$$W_n^* := \arg \min \frac{1}{n} \sum_{t=1}^n (Y_t - F_W(Z_t))^T \Gamma_0^{-1} (Y_t - F_W(Z_t))$$

then we have

$$\lim_{n \rightarrow \infty} \sqrt{n}(W_n^* - W^0) = \lim_{n \rightarrow \infty} \sqrt{n}(\hat{W}_n - W^0) = \mathcal{N}(0, I_0^{-1})$$

We remark that \hat{W}_n has the same asymptotic behavior than the estimator generalized least square estimator with the true covariance matrix Γ_0^{-1} which is asymptotically optimal (see for example Ljung [2]), so the proposed estimator is asymptotically optimal too.

4 Conclusion

In the linear multidimensional regression model the optimal estimator has an analytic solution (see Magnus and Neudecker [3]), so it doesn't make sense to consider minimization of a cost function. However, for the non-linear multidimensional regression model the ordinary least square estimator is sub-optimal if the covariance matrix of the noise is not the identity matrix. We can overcome this difficulty by using the cost function $U_n(W)$. The numerical computation and the empirical properties of these estimator have been studied in a previous article (see Rynkiewicz [4]). In this paper, we have given a proof of the optimality of the estimator associated with $U_n(W)$. This is then a good choice for the estimation of multidimensional non-linear regression model with multilayer perceptron.

References

- [1] Fukumizu, K., A regularity condition of the information matrix of a multilayer perceptron network, *Neural Networks*, Vol.9, 5:871–879, 1996
- [2] Ljung, L. *System identification : Theory for the user*, Prentice Hall, 1999
- [3] Magnus, J., Neudecker, H. *Matrix differential calculus with applications in statistics and econometrics* J. Wiley and Sons, New York, 1988.
- [4] Rynkiewicz, J., Estimation of Multidimensional Regression Model with Multilayer Perceptron. In J. Mira and A. Prieto, editors, proceedings of the 8th international workshop on artificial neural networks (IWANN 2003), Lecture Notes in Computer Science 2686, pages 310-317, Springer-Verlag, 2003.
- [5] Yao, J.F., On least square estimation for stable nonlinear AR processes, *The Annals of Institut of Mathematical Statistics* 52:316-331, 2000